

#38

EXHIBIT ONE

order to achieve the proper efficiency of translation. These rules are much less clear for viral RNAs, due to the interfering and overshadowing influences of other processes such as replication and packaging which require particular sequences.

The second manifestation of codon usage preference occurs in the selection of synonymous codons by different tRNA species charged with the same amino acid. Such a population of tRNAs is known as a family of "iso-accepting" tRNAs. As would be expected for strongly expressed genes, there is a clear correlation between the relative amounts of different iso-accepting tRNA species and the use of corresponding codons (Ikemura, 1981). In other words, the higher the concentration of a particular iso-accepting tRNA, the more often the corresponding codon appears in the sequence of the strongly expressed gene; on the other hand, this means that translation may be modulated and controlled by rare codons, for which the corresponding tRNAs occur in trace amounts only. Such codons may be AUA, coding for isoleucine, CUA (leucine), CGG and CGA (arginine), and GGA (glycine). These codons are marked by an arrow in Table 7-4, and, indeed, they are hardly used at all. It is not known whether organisms really use this mechanism to control gene expression.

It should be noted here that the choice of a codon for which there is a limited supply of a corresponding charged tRNA would inevitably cause an imbalance in the tRNA population of an organism. This in turn would not only slow down translation but would also make the system more prone to errors. One may imagine, for example, a competition between correct and false tRNAs at the ribosome A site prepared for the entry of an aminoacylated tRNA. If, due to its low concentration, the correct tRNA were too slow to interact, the false tRNA would associate with the ribosome and, hence, a false amino acid would be incorporated into the growing polypeptide chain. This process may even be associated with alterations of the reading frame if the structure of the

false tRNA prevents the proper entry of the next tRNA, and this has, indeed, been observed for several suppressor tRNAs. It is the basis of the phenomenon known as frameshift suppression. Suboptimal translation conditions of this kind have been artificially induced *in vivo* by starving bacterial cells for certain amino acids or *in vitro* by the addition of certain tRNAs to cell-free systems (Roth, 1981; Weiss und Gallant, 1983).

The significance of an appropriate codon choice for the expression of foreign genes in heterologous organisms has never been convincingly documented; nevertheless, especially since other unknown parameters may affect heterologous gene expression, the rules mentioned in this section should be followed as closely as possible in order to approach natural conditions. For chemically synthesised genes, for example, codons should be selected in accordance with the frequencies with which such codons occur in the desired host organism (cf. Section 11.2.2.1).

## 7.4 Construction of Expression Vectors

Several strategies using regulatory sequences discussed in the preceding sections have been pursued to optimise the expression of genes. In principle, these strategies are aimed at the construction of vectors allowing the synthesis either of fusion proteins comprising vector and insertion sequences (Fig. 7-44A) or of pure proteins exclusively encoded by the insertion (Fig. 7-44B). The first construction is referred to as a translational fusion, the second as a transcriptional fusion. The following selected examples will clarify this distinction.

### 7.4.1 Synthesis of Fusion Proteins

In order to obtain a hybrid protein, the foreign DNA must be inserted into an expressible vector gene in such a way that the reading frame in this

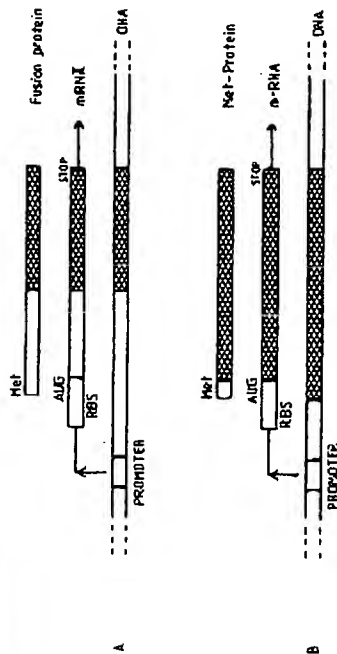


Fig. 7-44. Construction of expression vectors.

Two approaches are shown, namely the formation of fusion proteins (A), and the formation of native proteins (B) from recombinant DNA. RBS signifies a ribosomal binding site. Met-protein indicates that proteins obtained from recombinant DNA by approach (B) always carry an N-terminal methionine residue. Bacterial sequences are represented as open, eukaryotic sequences as hatched bars.

gene is conserved. The synthesis of hybrid mRNA is initiated by the prokaryotic promoter and its translation is controlled by the corresponding ribosome binding site. The first practical application of fusion proteins allowed the expression of rat insulin, rat growth hormone, and human growth hormone, and demonstrated for the first time that bacteria are, indeed, capable of expressing eukaryotic coding sequences.

#### 7.4.1.1 Expression of Rat Insulin

The starting point in this case was the insertion of a rat insulin cDNA into the *Pst*I site of pBR322 by homopolymeric poly(dG)-poly(dC) tailing (Villa-Komaroff *et al.*, 1978). The variable lengths of these tails guaranteed that at least one in three clones contained the right reading frame; however, since the cDNA could be inserted in two

Howard M. Goodman  
Boston

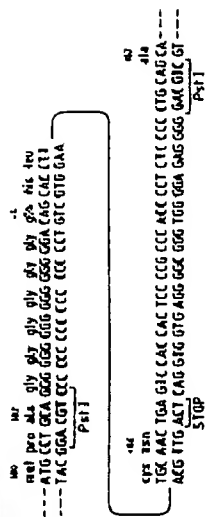


Fig. 7-45. Insertion of the rat insulin DNA sequence into the *Pst*I site of the  $\beta$ -lactamase gene of pBR322. The two *Pst*I sites, and amino acids 182 and 183 of  $\beta$ -lactamase, which are separated by the insertion, are printed in bold-face. The insulin insertion begins with amino acid Gln (position +4) in the B chain, and ends with aspartate of the proinsulin. The order of the insulin peptides is pre-B-C-A. (Villa-Komaroff *et al.*, 1978).

different orientations, only one sixth of the clones containing the desired insulin insertion would also make insulin. In spite of these obvious limitations, cloning by homopolymeric tails was the method of choice because the exact sequence of the cDNA was not known and the desired constructions therefore, could not be planned in advance (*cf.* also Section 3.2). The structure of one the rat insulin clones is shown in Fig. 7-45. Starting with position 182 (ala) the sequence of the  $\beta$ -lactamase gene then proceeds with polyglycine and eventually reaches the insulin sequence at amino acid "+4" (gln) of proinsulin. The desired fusion protein was detected by immunological techniques (see Section 11.2.3.2).

#### 7.4.1.2 Expression of Rat Growth Hormone and the Structural Protein VP1 of Foot and Mouth Disease Virus

A much more direct strategy was pursued for the construction of vectors coding for rat growth hormone (Seeburg *et al.*, 1978). The rat growth hormone cDNA possesses a single *Pst*I site at position "+24" of the prepeptide region, which allowed it to be annealed with the *Pst*I site of the  $\beta$ -lactamase gene of pBR322 in such a way that the reading frame was conserved (Fig. 7-46); in addition, the strategy employed for the construc-

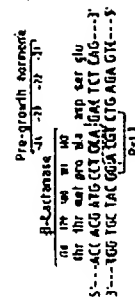


Fig. 7-46. DNA sequence in the vicinity of the *Pst*I site of a hybrid vector containing a fusion of the  $\beta$ -lactamase gene with the gene for rat growth hormone.

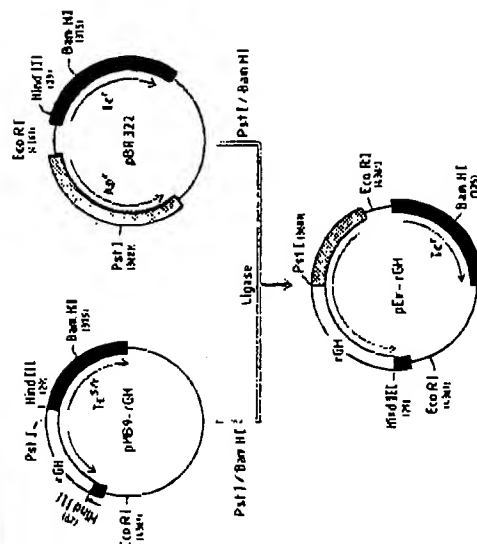


Fig. 7-47. Expression vector for rat growth hormone (rGH).

The expression plasmid pEx-rGH is constructed by replacing the small *Pst*I-BamHI fragment of plasmid pMB9-rGH by the smaller *Pst*I-BamHI DNA fragment of pBR322. Numbers in brackets refer to co-ordinates in plasmid pBR322. Arrows indicate the direction of transcription or translation. In pMB9-rGH, the promoter of the tetracycline resistance gene is interrupted by the rGH insertion; expression of tetracycline resistance is therefore markedly reduced. Transformants are resistant to 5 µg tetracycline/ml, while full expression in pBR322 or the expression vector pEx-rGH allows selections with 20 µg/ml. Tc<sup>r</sup> regions are represented as black, Ap<sup>r</sup> regions as hatched, and rGH regions as open bars. (Seeburg *et al.*, 1978).

395 amino acids in length. It comprises 181 N-terminal amino acids derived from the  $\beta$ -lactamase gene and 214 C-terminal amino acids of the pre-growth hormone and was, indeed, detected as a protein with a molecular weight of 46 000 in a mini-cell test system (see Section 11.2.3.3); however, the amount of  $\beta$ -lactamase produced by one-fifth the amount of hybrid protein was only pBR322.

Another example is the synthesis of a fusion protein containing a part of protein VP1 of Foot and Mouth Disease Virus (Fig. 7-48; Küpper *et al.*, 1981). In this case, cloning started with the insertion, into vector pLc24 cut with *Bam*HI and *Hind*III, of an 849 bp *Bam*HI-*Hind*III fragment coding for amino acids 9 to 292 of the desired protein (Fig. 7-28). The fusion protein obtained

was 395 amino acids in length and consisted of 98 N-terminal amino acids of MS2 replicase, 284 amino acids of the desired viral protein, and thirteen plasmid-derived amino acids added because of read-through into neighbouring vector sequences.

#### 7.4.1.3 Expression of Human Growth Hormone

Suitable restriction sites are rarely positioned such that they are located at the beginning of a structural gene and also allow this gene to be inserted into the vector gene in the correct reading frame. Quite frequently it is necessary to design special constructions in which linker molecules play an important role. The following

Fig. 7-48. Structure of an expression vector structural protein VP1 of Foot-and-Mouth Virus (FMDV).

The expression plasmid pPL-VP1 is derived from mid pLc24 which contains the N-terminal part of the VP1 structural protein under the control of the *P*<sub>L</sub> promoter. The structure of the FMDV cDNA and the pLc24 are shown. The VP1 structural protein is flanked by *Bam*HI and *Hind*III sites in the centre. Show bottom are the sequences around the *Bam*HI and *Hind*III sites in expression vector pPL-VP1. The *Hind*III site is at position 2105 of the pBR322 sequence.

example of a human growth hormone (hGH) expressed under the control of promoter may illustrate the point (Marria 1979). The starting material in this case was pPrpEDS-1 (Fig. 7-16) with a *Hind*III site the codon for amino acid 92 of the this *Hind*III site had been joined with a fragment flanked by a *Hind*III site in untranslated region of the cDNA for growth hormone, the correct reading frame have been lost (Fig. 7-49). It was th necessary to manipulate the *Hind*III pPrpEDS-1 to shift the reading frame by a in the recombinant molecule. This was plished by filling in the 5' protruding en Klenow fragment of *E. coli* DNA polym and adding a synthetic DNA decamer contained a *Hind*III site. As shown in Fig cleavage of the new plasmid pPrpEDS *Hind*III and subsequent ligation with *Hind* bGH cDNA conserved the correct reading

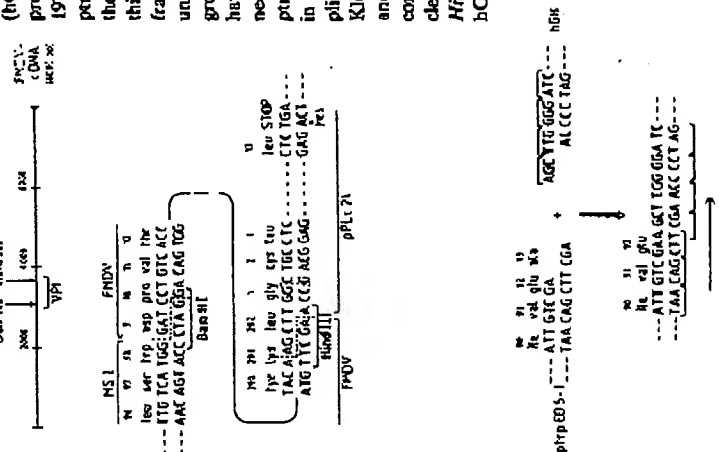


Fig. 7-49. Fusion of a *Hind*III site in plasmid pPrpEDS-1 with a *Hind*III-flanked cDNA fragment of human growth hormone (hGH). Brackets indicate the reading frames, the arrow denotes the direction of translation. The bGH section derived from the 5' untranslated region.



by partially digesting  $\lambda$ phac3 DNA with *Hae* III. The resulting mixture of DNA fragments contained a 203 bp *Hae* III fragment containing the entire *lac* control region and the first seven amino acids of  $\beta$ -galactosidase (cf. Fig. 7-7). The mixture of DNA fragments was then ligated with pBR322 DNA which had been linearized with *Eco* RI and subsequently filled-in to convert its protruding 3' ends to flush ends. Ligation of filled-in *Eco* RI ends with blunt *Hae* III ends generated new *Eco* RI ends in the recombinant molecules at the sites of fusion (see also Section 2.1.2.1).

Transformants obtained from this DNA mixture containing the desired *Hae* III fragment were identified as blue colonies on agar plates containing Xgal (Fig. 7-6). This screen did not distinguish between the two possible orientations of the *Hae* III fragment; however, since there was an asymmetrically positioned *Hha* I site directly following the stop codon of the *lacI* gene on the *Hae* III fragment (Fig. 7-7), it was easy to determine the orientation of the inserted *Hae* III fragment. The desired orientation was found in vector pBH10 (Fig. 7-52), in which *lac* transcription proceeds toward the tetracycline resistance region.

An unusual procedure was used to selectively remove the distal *Eco* RI site in pBH10. *E. coli* RNA polymerase binds to promoter regions in the absence of nucleoside triphosphates. In pBH10 binding occurs at the *lac* promoter and

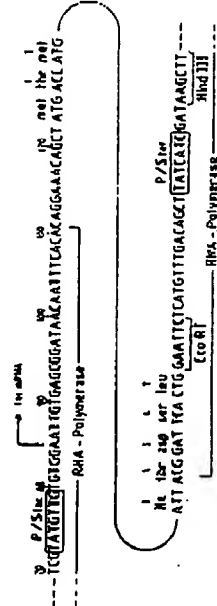
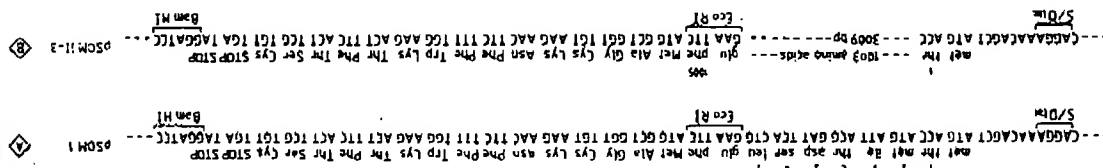


Fig. 7-43. DNA sequences between the *lac* and *trp* promoter regions in plasmid pBH10. The Pribnow-Schaller boxes are framed, the proximal *Eco* RI site of the *trp* promoter and the *Hha* III site are bracketed. Binding sites for RNA polymerase, which extend approximately 35 bp to the left and right of the Pribnow-Schaller boxes, are indicated by brackets. It is apparent that the proximal *Eco* RI site lies in a region protected by RNA polymerase. For the numbering in the *lac* region see legend to Fig. 7-7.

this procedure should have yielded functional somatostatin; nevertheless, all attempts to do the hormone in various bacterial extracts were unsuccessful. These negative results probably resulted from proteolytic cleavage of the hormone (cf. Section 7.5).

A new plasmid, pSomII, was constructed to hope that the presence of a large peptide would prevent this proteolytic attack (Fig. 7-52). In plasmid the smaller *Eco* RI-*Pst* I fragment the *lac* region of pSomI was replaced by the corresponding *Eco* RI-*Pst* I fragment of pBR. Transformants were selected for ampicillin resistance and screened on Xgal plates for the absence of *lac* operator DNA. A *lac* region containing the entire control region and the codons for 100 of 1 021 amino acids of  $\beta$ -galactosidase, was to replace the missing *lac* region. The *Eco* RI fragment of 7.45 kb was obtained by  $\lambda$ phac3 DNA. As shown in Fig. 7-54B, the reading frame was retained in this construct and a fusion protein of 1 020 amino acids with amino acid sequence of somatostatin at C-terminal end was obtained. When total cell proteins of saltably induced bacteria were treated with cyanogen bromide, somatostatin activity was detectable. The yield in uninduced was estimated to be on the order of 0.001-0.01 of the total protein. This low yield reflects the basal level of transcription from a fully repressed *lac* promoter. Induction with IPTG led to a sevenfold increase in somatostatin yields. Induction experiment confirmed the sequence data indicating that in pSomII synthesis of somatostatin was regulated by the control sequences. However, the induction was approximately tenfold lower than had

Fig. 7-54. Nucleotide sequence of *lac* fusion genes. (A) shows a fusion with only seven N-terminal acids of  $\beta$ -galactosidase, (B) a fusion with 1000 acids of  $\beta$ -galactosidase. Both fusions contained additional amino acids derived from the *Eco* I (Itakura et al., 1977).



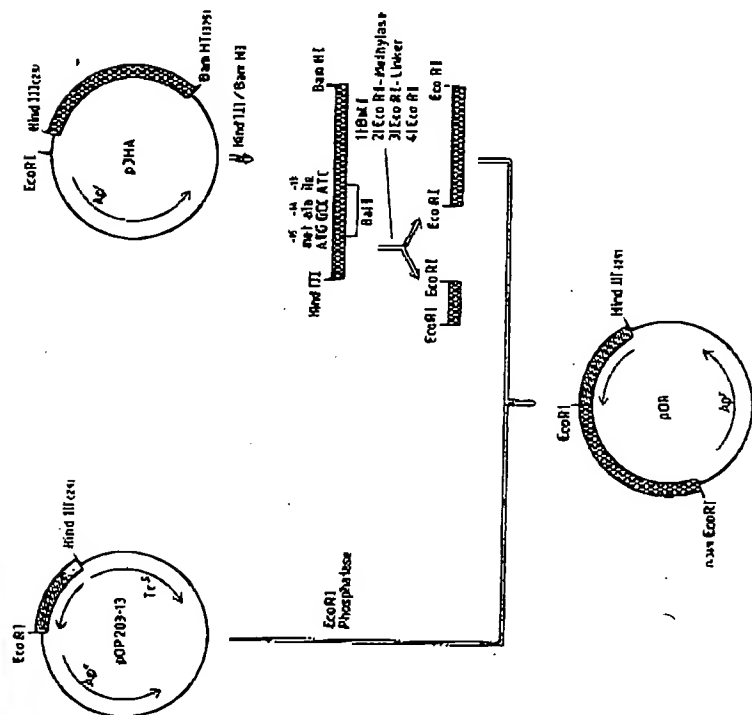


Fig. 7.55. Construction of a plasmid for the expression of the haemagglutinin (HA) gene of human influenza virus type A/Japan/305/57 (subtype H2). The modifications of the HA DNA fragment which is obtained from pHHA by Hind III-Bam HI digestion are described in detail in the text. Lac control regions are indicated by a stippled bar. HA DNA sequences by a cross-hatched bar (Heiland and Gething, 1981).

expected. Similar observations subsequently have been made with other expression plasmids based on *lac* control elements. There are several possible explanations for this phenomenon, including the selective cleavage of the foreign protein by bacterial proteases, insufficient solubility of the fusion protein during cyanogen bromide cleavage and the instability of the recombinant plasmid.

#### 7.4.1.5 Construction of Expression Plasmids for Influenza Virus Specific Sequences

This case deals with the expression of a DNA copy of an RNA fragment coding for the haemagglutinin (HA) protein of human influenza virus strain A/Japan/305/57 (subtype H2) (Heiland and

Gething, 1981). The vector used was plasmid pOP203-13 (Fig. 7-9) which contains, between the *Eco* RI and *Hind* III sites of pBR322, the same 203 bp of the *lac* control region as pH10 (Fig. 7-52). The direction of *lac* transcription is anticlockwise, i.e., in the direction of the  $\beta$ -lactamase gene (Fulter, 1982), which means that any DNA inserted into the single *Eco* RI site of this plasmid will be controlled by the *lac* promoter.

The haemagglutinin gene to be expressed was inserted between the *Hind* III and *Bam* HI sites of pBR322 in plasmid pHHA (Fig. 7-55). It codes for the entire 560 amino acids of the haemagglutinin protein and eleven nucleotides of the 5' untranslated region. This sequence must be modified before it can be inserted into the *Eco* RI site of pOP203-13. A *Bal* I site comprising the ATG start codon of the HA gene is important. The DNA fragment obtained by *Hind* III and *Bam* HI digestion is further cleaved by *Bal* I treatment to yield two sub-fragments. The mixture of fragments is first treated with *Eco* RI methylase in order to methylate internal *Eco* RI sites and to render them resistant to *Eco* RI digestion. *Eco* RI linkers are then added by ligation. Following *Eco* RI digestion the sub-fragments are separated from each other and the larger fragment is inserted into the *Eco* RI site of the expression vector pOP203-13 (Fig. 7-55).

Cloning of the large *Bal* I fragment was expected to yield a fusion protein with the structure shown in Fig. 7-56, containing seven N-terminal amino acids derived from  $\beta$ -galactosidase, three amino acids coded for by the linker, and 560 amino acids of the haemagglutinin gene. Two N-terminal amino acids of the leader sequence of the HA gene were removed by this cloning procedure. Again, the two possible orientations for the inserted gene could be easily distinguished by suitable digestions. Three of the clones obtained expressed antigenic determinants of haemagglutinin, as shown by solid phase radioimmunoassay. The nucleotide sequences of all three clones confirmed that they preserved the correct reading frame, but also showed that they did not have the

expected sequence of the hypothetical pOR (Fig. 7-56) at the site of fusion. For unl reasons fifteen amino acids of the signal p and the first ten to fifteen amino acids mature protein were missing. Perhaps the e otic hydrophobic signal sequences were g erated by the *E. coli* host organism. i particular case. By way of contrast, other l phobic signal sequences, such as that of l preproinsulin, have been found to be quite in *E. coli* (Chan *et al.*, 1981).

#### 7.4.1.6 General Technique for the Construction of Expression Vectors for Fusion Proteins

In the examples discussed so far, a restriction site within the region of the ATG start codon was always positioned in such a way that the HA gene could be inserted into the *lacZ* gene (or another suitable gene) either conserved the correct reading frame. In the case of the somatostatin gene this was second case by suitably planning the chemical synthesis while it was mere coincidence in the case HA gene. In most cases, however, a convenient restriction site will not be available; the following procedure is therefore recommended for cloning and expression: the DNA to be expressed as a cDNA, is cut out and isolated in a parent plasmid. The example shown in Fig. 7-56 uses a *Pst* I digestion. The next step is to choose a suitable restriction site in the vicinity of the start codon. The DNA is first treated with a combination of the enzymes *Eco* III and *S* I (cf. also Fig. 2.1-9). Digestion conditions depend on the distance between the start codon and the *Pst* I site (in our example) and the start codon, and must be determined for every individual case. In our example, the *Pst* I site is then added to the fragment, so that it can be inserted into the *Eco* RI site of a suitable vector. Before the fragment containing the HA gene is cloned, it is digested with *Eco* RI and a restriction enzyme which cleaves at site X

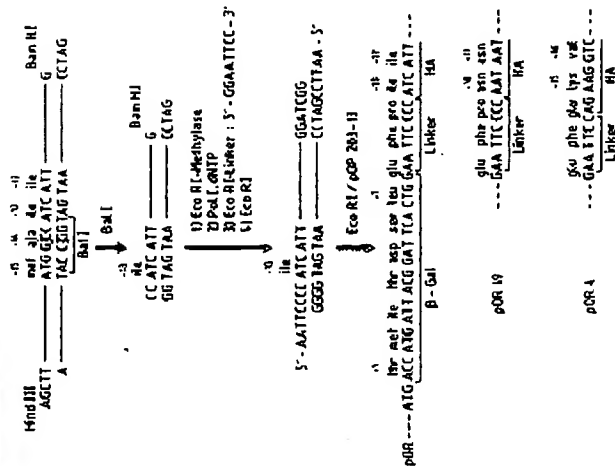


Fig. 7-56. Construction of an expression plasmid by linker technology.

A Hind III-Bam HI fragment of plasmid pHA (Fig. 7-55) is cleaved into two fragments by digestion with Bam HI. Only the flanking sequences of the larger of the two fragments are shown. Following an Eco RI methylation treatment, the protruding Bam HI site is filled-in in order to allow the subsequent addition of Eco RI linkers. A further Eco RI digestion only attacks the Eco RI sites within the linker, but not the internal methylated Eco RI site. This DNA fragment is then ligated into the Eco RI site of plasmid pOP203-13 (Fig. 7-9). Shown is the expected structure of the fusion protein (pOR) consisting of seven amino acids of  $\beta$ -galactosidase, three amino acids encoded by the linker, and two amino acids out of a total of 360 from the haemagglutinin. The actual experiment did not yield clones with the expected structure; instead, plasmids were obtained, which begin with sequences of the mature HA protein and which do not contain the hydrophobic leader sequence. Part of the structures of two of these plasmids, pOR19 and pOR4, are also shown.

the gene to be cloned. This yields a defined right-hand molecular end which can be used at a later stage to reconstruct the entire gene. Although the left end of the fragment is defined by an Eco RI site, the distance between this site and the ATG start codon varies in different molecules. Ligation with a suitable vector will therefore yield a wide range of different clones with varying distances between start codon and Eco RI site. In addition, the insertions may not be

in the correct reading frame. Those clones in the mixture which contain the correct reading frame, can be identified by exploiting the phenomenon of  $\alpha$ -complementation (cf. Section 2.4.2.3). As in the case of M13 cloning (Section 2.4.2), a number of plasmids, known as pUC plasmids, have been developed for this purpose (Fig. 7-58). These plasmids contain the lac regulatory region and a part of the lacZ gene which codes for the 59 N-terminal amino acids of  $\beta$ -galactosidase (Vieira

and Messing, 1982). The corresponding host strain (JM83) carries the deletion M15 of the lac operon, which removes amino acids 11-41 of  $\beta$ -galactosidase, but retains the entire C-terminal part of the enzyme. Each incomplete lacZ gene will direct the synthesis of an inactive polypeptide. Together, these polypeptides will be capable of complementing each other by forming aggregates. The resulting enzymatic activity can be detected on Xgal indicator plates as described above (Fig. 7-6).

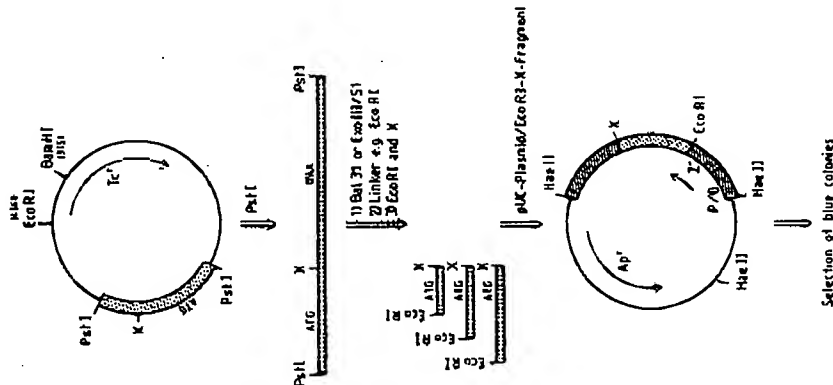


Fig. 7-57. General approach for the construction of expression vectors directing the synthesis of proteins.

The starting material can be a cDNA clone. Insertion has been removed from the plasmid (b digestion in this example). By treatment with a restriction enzyme (X) a combined EcoRI/XbaI site is positioned close to the start of the cDNA. The example of EcoRI/XbaI is shown here may illustrate the use of EcoRI/XbaI. The next step contains an internal EcoRI site. The next step is EcoRI digestion followed by digestion with a restriction enzyme (X) which should preferentially cleave the cDNA. The mixture of DNA fragments is then ligated into a pUC vector (cf. also Fig. 7-58). Of course, the cleavage site of EcoRI/XbaI is present within the polylinker of the pUC plasmid. A wide spectrum of pUC plasmids is available, not being difficult to find a suitable vector. Once a suitable clone is found, the lacZ gene can be used for the insertion of the portion of the gene in question, which can be obtained from the original cDNA clone. The lacZ gene and regions are represented by hatched bars; the lacZ gene is represented by a stippled bar.



## 292 7 Expression Vectors in Prokaryotes



pUC plasmids (cf. Fig. 2, 4-22) are derived from the 1297 bp PvuII-Eco RI fragment of pBR322, which contains the origin of DNA replication (*ori*) and the coding region of  $\beta$ -lactamase (*Ap<sup>r</sup>*). *Pst*I, *Hind*III and *Acl*I restriction sites were removed by mutagenesis; pUC plasmids carry a 433 bp *Hae*II fragment with *lac* control elements (*lac* promoter (P) and operator (O); open bars) inserted into the *Hae* II site in the immediate vicinity of the replication origin at position 2352; in addition, they contain the coding region for a functional  $\beta$ -galactosidase  $\alpha$ -peptide (*lacZ*) (hatched bar). Short polynucleotide regions within this region provide multiple recognition sites for various restriction endonucleases. Amino acids encoded by polynucleotide insertions are printed in Italics. Numbers in parentheses are pUC18 co-ordinates (Appendix B-4; Weiss and Messing, 1982; Yanisch-Perron *et al.*, 1985).

In summary, the successful synthesis of hybrid proteins with prokaryotic and eukaryotic components has been described in detail for several systems. This approach has considerable advantages, particularly the fact that fusion proteins with the large, 1 000 amino acids-long *N*-terminus from  $\beta$ -galactosidase often are insoluble within the bacterial cell. Such fusion proteins thus are protected from proteolytic degradation (see below) and are easily purified; however, it should be kept in mind that this strategy also has its limitations. There is no doubt that it permits the detection of antigenic determinants in the fusion protein. Actually, it allowed the initial demonstration of the possibility of expressing eukaryotic DNA sequences in prokaryotes; however, if the eukaryotic proteins in question are to be obtained in a pure form, the original protein must be separable from the bacterial component of the chimeric fusion products. Cyanogen bromide cleavage, which was used in the case of somatostatin, is restricted to proteins, such as some proinsulins, which do not contain internal methionine residues. In other cases enzymatic cleavage must be employed. Since the codons for suitable specific amino acids (for example arg and lys residues for tryptic cleavage) usually are not found in desired positions on the vector, this

## Bacteria

In contrast to procedures described above, this process, known as transcriptional fusion, directly affects the production of unique, non-polyepitopes. In this case protein synthesis does not initiate from the first methionine residue of the prokaryotic leader peptide, such as *lacZ* or *lacY*, but from the first methionine of the diphthery toxin polypeptide itself (Figs. 7-44B). Biologically

#### 7.4.2.1 The lac System

M. Prashne and co-workers have developed a concept of a portable *lac* promoter which can be placed at a suitable distance in front of a de-

structural gene to allow the synthesis of the corresponding gene product as a pure unfused protein. The *lac* operon contains a suitable *Alu*I site between the ribosomal binding site and the start codon of the *lacZ* gene. A 95 bp long *Alu*I fragment containing almost the entire *lac* promoter region, the initiation site for mRNA synthesis, the S/D sequence, and five additional base pairs therefore can be isolated from chromosomal DNA (Fig. 7-7). This fragment also can be obtained from plasmid pGL101, in which this fragment is flanked by an *Eco*RI and a *Pvu*II site

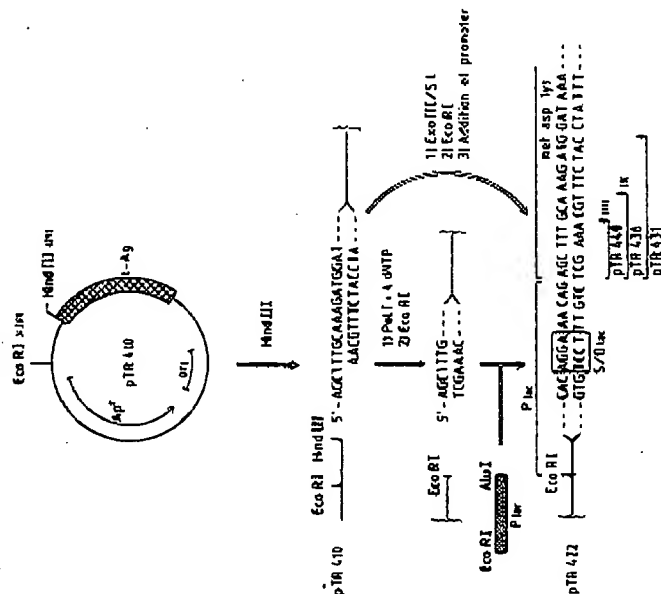


Fig. 7-59. Expression plasmid for SV40 t-antigen. Vector pTR410 consists of parts of pBR322 comprising the coding region for  $\beta$ -lactamase (*Ap*<sup>r</sup>) and the origin of DNA replication; in addition, it contains SV40 sequences with the entire coding region for the small t-antigen (cross-hatched bar). Manipulations near a *Hind*III site upstream of the start codon for t-antigen are described in detail in the text. *Pvu*II signifies the portable *lac* promoter fragment (stippled bar; see also text). Small numbers at the arrow heads of deletions pTR440 and pTR436 indicate the distances between the S/Diac regions and the start codon.

(Fig. 7-10) (Lauer *et al.*, 1981). Since *Pvu*II (CAG/CTG) recognises a hexanucleotide sequence comprising the *Alu*I recognition sequence (AGACT) and produces blunt ends at the same position as *Alu*I, *Eco*RI/*Pvu*II cleavage of pGL101 yields a DNA fragment with the desired blunt *Alu*I ends immediately downstream from the *lac* S/D sequence. This fragment contains an S/D sequence, but lacks an ATG codon, and must be placed at the proper distance upstream from a structural gene. For this purpose, the gene in question should preferably contain a unique

recognition sequence immediately 5' of its ATG codon, which is, of course, rare. In the case of small t-antigen of SV40, however, the *Hind*III site in the SV40 *Hind*III-B fragment (which is 1169 bp in length), is only twelve bp upstream of the ATG start codon of t-antigen. When this DNA fragment with filled-in *Hind*III ends was annealed with a portable promoter, the construction designated pTR422, shown in Fig. 7-59, was obtained (Roberts *et al.*, 1979b).

In order to obtain clones in which the distance between the ATG codon and the *lac* S/D sequence is shortened, the DNA first was digested with *Hind*III and then subjected to a partial digestion with *S*1 nuclease treatment and the molecule was circularised after *Eco*RI cleavage and addition of the portable promoter. This procedure yielded

inserts of variable length due to the unспецифичность of the exonuclease III reaction. The clones then screened for expression of the protein. In this example, clone pTR436, in the distance between the S/D sequence and the ATG start codon was a *sp.* was particularly (Fig. 7-59). Plasmid pTR440 is only weakly and the starting plasmid pTR422 and a deletion pTR431, were completely inactive.

In a similar case the starting materials were same as those described above, namely portable *lac*/UV5 promoter, and the *Hind*III-B fragment with the coding region of the small t-antigen (Thummler *et al.*, 1981) strategy employed, however, differed from described above in that a *Hind*III link introduced between the S/D sequence and the ATG start codon (Fig. 7-60). This *Hind*III

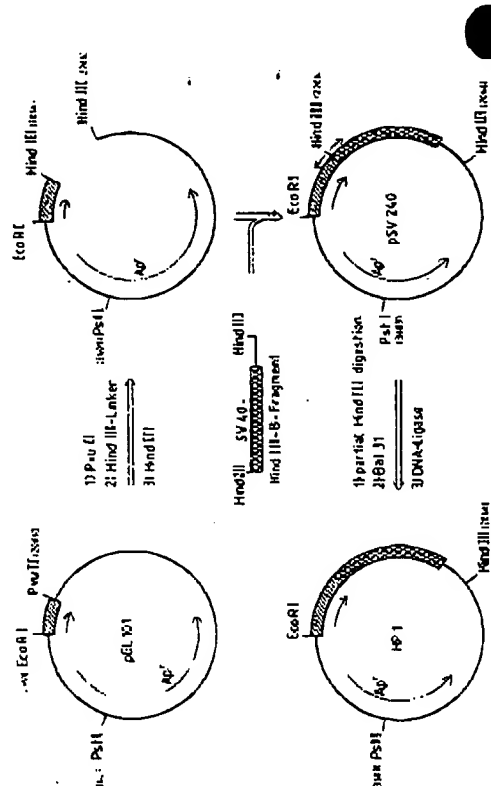
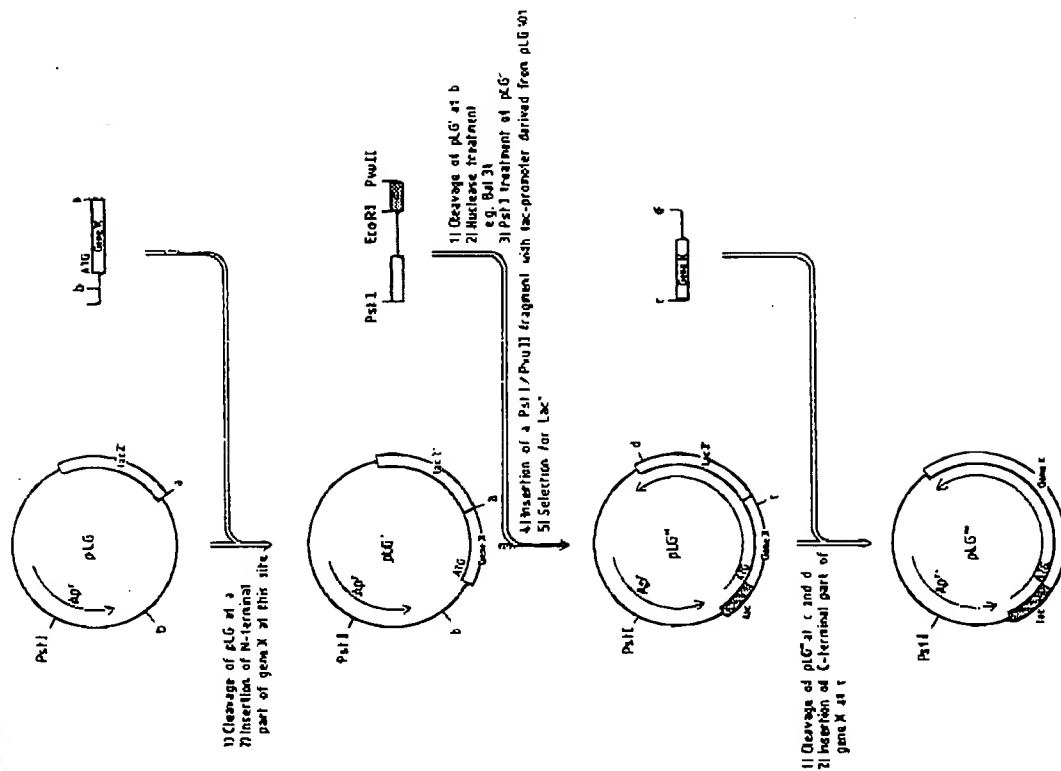


Fig. 7-60. Construction of an expression vector for SV40 t-antigen. The *Pvu*II site of pGL101 (Fig. 7-10) is converted into a *Hind*III site by using *Hind*III linkers, which allow subsequent insertion of the SV40 DNA fragment. In contrast to the construction shown in Fig. 7-59, the *lac* promoter (cross-hatched bar) is not inserted into the SV40 DNA fragment (stippled bar), as well as the S/D sequence within the portable *lac* promoter (stippled bar) are shortened by exonuclease treatment (indicated by arrows extending from the *Hind*III sites). Numbers in brackets are pBR322 co-ordinates. Directions of transcription are indicated by arrow heads of the plasmid circles.





$\beta$ -Galactosidase  
SV40T-Ag (pTR30):  
---TAA(AATT)CACA(CAGGAACAGCAATG)---  
---TAA(AATT)CACA(CAGGAACAGCAATG)---  
 $\beta$ -Globin (pLG 102-7):  
---TAA(AATT)CACA(CAGGAACAGCAATG)---  
pEF1 (pLG 101):  
---TAA(AATT)CACA(CAGGAACAGCAATG)---  
pIF (pLG 101):  
---TAA(AATT)CACA(CAGGAACAGCAATG)---  
SV40T

Fig. 7-65. Ribosomal binding sites in different expression vectors.

All structures contain the same SD sequence also found in the lac system (boxed). The point of translation is indicated by a vertical line (Guarente *et al.*, 1980a).

4- Fig. 7-66. General approach for constructing optimally expressing clones in the lac system. The 5' terminal part of a gene X to be expressed is introduced into a restriction site "a" of a plasmid containing the 3' terminal portion of the lacZ gene. Plasmid pLG<sup>+</sup> is then opened at "b", modified by nuclease treatment and ligated with a lac promoter fragment (obtained, for example, from pGL101, Fig. 7-10) which is flanked by PstI and PvuII sites. Maximally expressing clones are identified by selection for Lac<sup>+</sup> in a growth medium containing Xgal. The lacZ<sup>+</sup> region in these clones is then replaced by the 3' terminal portions of X by using site "c", which restores gene X (Guarente *et al.*, 1980).



Benno Müller-Hill,  
Cologne, FRG

Table 7-5. Expression yields of IGF-1/lacZ fusion proteins from different mutated plasmid constructions.

Plasmid	2	3	4	5	6	Cys	Leu	Thr	ACC	CTG	TGC	$\beta$ -gal activity units/cell (JM83)	ng/10 <sup>7</sup> cells SMC (HB101)
original sequence pUCmuSMCAori	Pro	Glu	Thr	Leu	Cys							0.4	1.4
blue colonies													
pUCmuSMCA	1	2	3	4	5	6	7	8	9	10			
1	CCC	GAA	ACT	CTG	TGT							3.1	33
2	CCT	GAA	ACT	TTC	TGC							2.6	45
3	CCA	GAG	ACG	TTC	TGC							0.9	35
4	CCA	GAG	ACG	TTC	TGC							0.9	40
5	CCT	GAA	ACT	TTC	TGT							2.9	33
6	CCT	GAG	ACG	TTC	TGT							1.2	58
7	CCG	GAA	ACG	TTC	TGT							1.9	50
8	CCG	GAA	ACA	TTC	TGT							1.2	65
9	CCA	GAA	ACG	TTC	TGT							1.1	32
10	CCT	GAG	ACT	CTA	TGT							2.3	22
white colonies													
pUCmuSMCA	11	12	13	14									
11	CCC	GAA	ACC	CTC	TGT							<0.1	0.10
12	OCT	GAA	ACC	CTC	TGT							<0.1	0.11
13	CCG	GAA	ACC	CTC	TGT							<0.1	0.10
14	CCA	GAA	ACC	CTC	TGT							<0.1	0.09

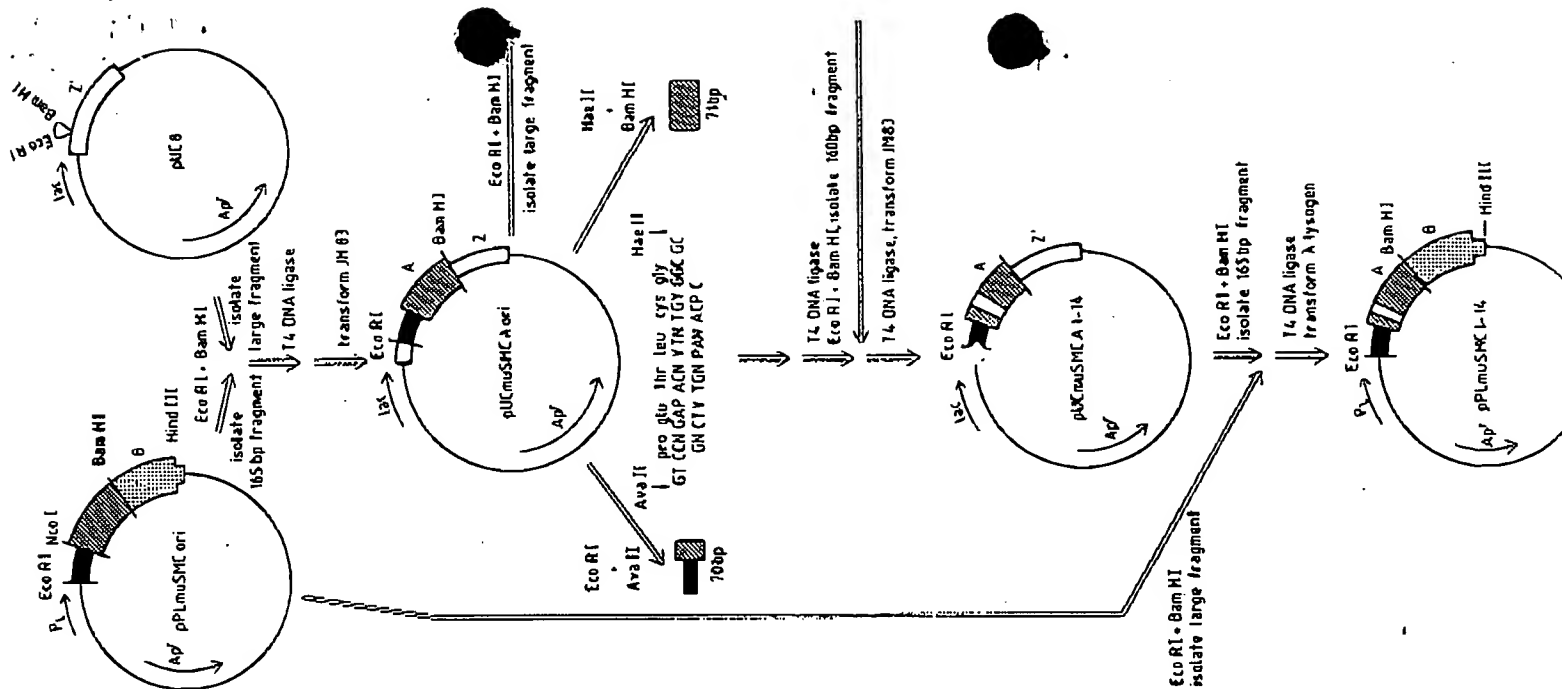
The table relates to plasmid constructions described in Figs. 7-66 and 7-67. SMC stands for somatomedin-C, a 70 amino acid protein found in human serum, also known as insulin-like growth factor I (IGF-I). Numbers 1 to 14 refer to fourteen plasmid colonies, ten of which displayed a blue and four of which displayed a white phenotype following plating on JM83 cells in the presence of Xgal plus ampicillin. x indicates the positions of mutations introduced in codons 2 to 6 of the IGF-I gene (Buell *et al.*, 1985).

tures of the respective ribosome binding sites are shown in Fig. 7-65. At equilibrium the corresponding bacteria synthesise between 5 000 and 15 000 molecules of the desired protein per cell. In individual cases the yields may be lower since the various proteins may differ in their stability within the host bacteria (Guarente *et al.*, 1980a) (see Section 7.5).

This concept can also be applied to the pUC family of vectors. In order to optimise expression of the IGF-I gene in *E. coli*, Buell *et al.* (1985) inserted a 165 bp fragment containing a ribosome binding site and the coding region of the first 33 amino acids of IGF-I into the 65 bp polylinker region of pUC8 (Fig. 7-67). Expression in this construction was under the control of the *lac* promoter, while translation could initiate either at the *lacZ* gene or at the IGF-I gene start codon. However, since translation from the *lacZ* AUG would quickly encounter a stop codon (Fig. 7-67),

Fig. 7-66. Construction of IGF-1/lacZ fusion vectors for improved expression of the IGF-I protein.

Vector pPLmuSMCAori contains a synthetic IGF-I gene (parts "A" and "B", hatched and stippled bars) preceded by a 66 bp fragment derived from bacteriophage mu (see Fig. 7-67) which provides the SD sequence (black bar). The construction is driven by the  $\lambda P_L$  promoter, but results in only low level expression of the desired protein. To improve expression, the N-terminal part "A" of the IGF-I gene is cloned into pUC8 to yield pUCmuSMCAori. An *Ava*II-*Hae*II fragment, (G/G(AorT)CC) (PuGCCG/Py), coding for amino acids 2-8 of IGF-I, is replaced by a synthetic mixture of DNA fragments containing all possible base substitutions which retain the amino acid sequence. The *Eco*RI-*Bam*HI fragments from plasmids isolated from blue colonies were isolated and reconstructed into pPLmuSMCA-14, as indicated. The open bar sector within part "A" of the IGF-I gene in plasmids pUCmuSMCA-14 and pPLmuSMCA-14 represents the synthetic fragment. N = one of the four possible bases, P = purines, and Y = pyrimidines (Buell *et al.*, 1985).



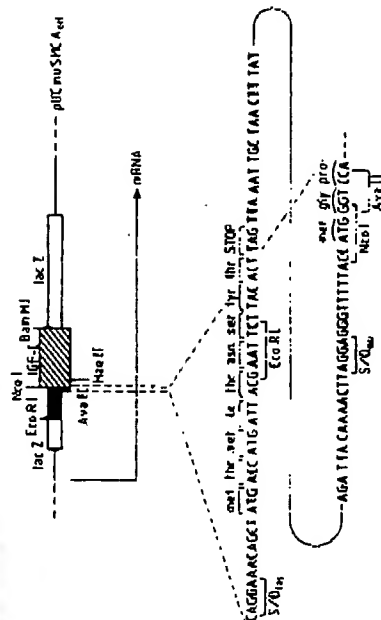


Fig. 7-67. Construction of IGF-1/lacZ gene fusions. The figure shows a section of plasmid pUC19 with an insertion between its Eco RI-Bam HI polylinker sites within the lacZ  $\alpha$ -peptide (open bar). The insertion comprises a 66 bp Eco RI-Nco I fragment containing the S/D sequence from the *nerf* gene of bacteriophage mu (Gray *et al.*, 1984) (black bar), and a 98 bp long Nco I-Bam HI fragment with the first half (part "A" in Fig. 7-66) of the coding region of the human IGF-1 gene (hatched bar). Transcription in this construction starts at the lacZ promoter present in the pUC19 portion of the vector, and covers the lacZ portion and the insert. Translation from this message can initiate both at the lacZ ribosomal binding site (S/D<sub>lac</sub>) and the normal binding site (S/D<sub>IGF</sub>) of the insert derived from bacteriophage mu. Ribosomes initiating at the lacZ ribosomal binding site will soon encounter a stop codon, while the second translation product continues through the IGF-1 coding region into the lacZ portion of pUC19. This construct yields only white plaques on *E. coli* strain JM83. Mutations introduced into the coding region within an Aval-Hind III fragment (*cf.* Fig. 7-66) result in a blue phenotype, indicating expression of the IGF-1/lacZ fusion peptide (Bucell *et al.*, 1985).

the only protein formed was derived from a fusion between the IGF-1 portion and the distal lacZ gene region (Fig. 7-67). Transfection of a construction containing genuine IGF-1 sequences into JM83 yielded only white colonies, indicating little or no  $\beta$ -galactosidase activity. In order to increase expression, a large number of mutants were generated by synthesizing a mixture of oligonucleotides that included all the 256 possible sequences encoding amino acids 2-6 of IGF-1 (Fig. 7-66). After re-insertion into the proper position in the fusion and transformation of *E. coli* strain JM83, approximately 500 out of 5 000 colonies were pale blue. The best of these, after reconstruction of the whole IGF-1 gene, produced more than 20 times more IGF-1 than did the wild-type construction (Table 7-5). These mutations affected a secondary stem structure around a ribosome binding site region and the increased

$\beta$ -galactosidase activity in the mutants thus confirms some of the conclusions mentioned in Section 7-2.

#### 7.4.2.2 The *trp* System

As in the lac system, the regulatory sequences of the *trp* operon can be used to create hybrid ribosome binding sites. In the *trp* system, the site of transcription initiation and the start codon of the *trpE* protein are separated by 162 bp known as the leader sequence (Fig. 7-14). This region codes for a peptide which is fourteen amino acids in length and plays a decisive role in the control of the *trp* operon. A *Taq* I site is situated between the corresponding S/D sequence and the ATG codon, allowing both parts of the ribosome binding site to be separated from each other. The ATG of the leader peptide can therefore be

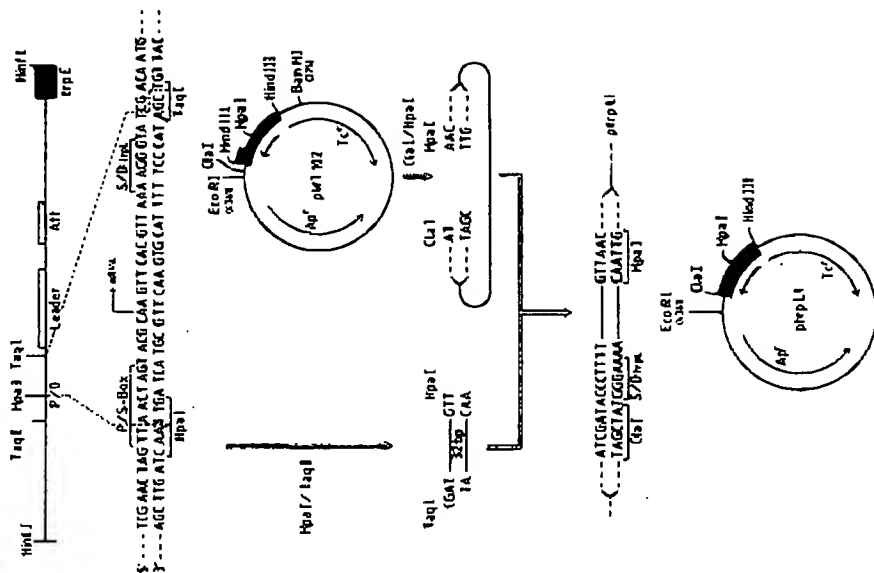


Fig. 7-68. Construction of an expression vector with *trp* regulatory sequences.

The top shows a portion of a Hind III DNA fragment with *trp* regulatory sequences (*cf.* also Fig. 7-15) between a Hpa I and a *Taq* I site containing the Pribnow-Schaller box (P/S) and the S/D sequence of the *trp* peptide (S/D<sub>trp</sub>). This 32 bp Hpa I-*Taq* I fragment is inserted into vector pMT102 opened with Hpa I and pMT102 and pMT101 (Fig. 7-18) are identical, but contain the *trp* insertion in opposite directions (arrows) unexpectedly, pMT102 shows a considerable tetracycline resistance, which is presumably due to the presence of a cryptic promoter in the *trp* Hind III fragment. By a fortunate coincidence, the *Taq* I site of the insertion contains a pair, and this regenerates the Cla I site. Vector pMT101 lacks the *trp* coding regions, and the Cla I site can there used directly for cloning of foreign DNA (Edman *et al.*, 1981).

replaced by the ATG of a eukaryotic gene.

A suitable plasmid which contains the entire *trp* promoter region, the *SD* sequence for the leader peptide, but no other parts of the *trp* operon was constructed from plasmid pWT102 (Fig. 7-18) (Edman *et al.*, 1981). Digestion with a combination of *Hpa*I (GTT/AAC) and *Taq*I (T/CGA) yielded a 32 bp fragment which contained the desired *SD* sequence and the initiation site for transcription (Fig. 7-68). In a parallel experiment all *trp*-specific regions upstream of the *Hpa*I site were removed from the same plasmid, pWT102, by digestion with *Hpa*I and *Cla*I and replaced by the short *Hpa*I-*Taq*I fragment. This construction presented no problems since the tetrameric *Taq*I recognition sequence (T/CGA) is part of the hexameric recognition sequence of *Cla*I (AT/CGAT), and since both enzymes cleave in the same pattern with protruding 5'-CG termini. Since there was an AT pair in the immediate vicinity of the original *Taq*I site, the *Cla*I site was regenerated. The resulting plasmid, pTP-L1, was opened at its unique *Cla*I site to allow insertion of a foreign gene in the immediate vicinity of the *SD* sequence. The *Cla*I site is particularly well suited for this purpose since it can accommodate various DNA fragments with protruding 5'-CG ends, such as those obtained by *Hpa*II (C/CGG), *Taq*I (T/CGA) and *Acl*I (GT/CGAC) cleavage. This strategy was tested and employed for cloning and expressing hepatitis B virus core antigen (HBcAg). The gene for this protein, which is 183 amino acids in length, was obtained from a suitable plasmid by *Hpa*I cleavage (Fig. 7-69). The start codon for HBcAg on a *Hpa*I fragment of 1005 bp was 15 bp away from one of the molecular ends of this fragment. The distance between the ATG and the *SD* sequence would have been too long, and therefore the usual modifications were carried out as described above. The DNA was first treated with exonuclease III to remove approximately ten base pairs and then made blunt-ended with S1 nuclease before *Bam*HI linkers were added. The commercially available decameric linkers do not only contain a

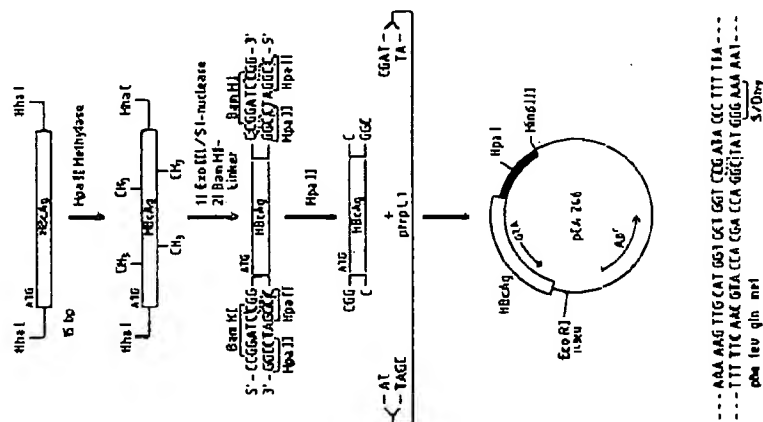


Fig. 7-69. Application of expression vector pTP-L1. A *Hpa*I DNA fragment containing the coding sequence for HBcAg was modified by a combined exonuclease III/S1 nuclease treatment. Following addition of *Bam*HI linkers, the construct was digested with *Hpa*I in preparation for cloning into the *Cla*I site of pTP-L1. The sequence around the ribosomal binding site at the bottom shows that the distance between *SD* sequence and ATG codon is 16 bp. The direction of transcription of the HBcAg sequences is indicated by an arrow.

*Bam*HI site but also two *Hpa*II sites. A subsequent *Hpa*II digestion completely removes the *Bam*HI recognition site and creates 5' overhanging ends which are compatible with *Cla*I ends. Since the tetrameric *Hpa*II recognition site occurs

Fig. 7-70. Physical maps of pKN402 and its derivatives. Plasmid pKN402 (shown on top in a linear presentation) is a mini-derivative of a temperature-sensitive replication mutant of plasmid R1d19. Plasmids pCP3 and pCP39 are derivatives of pKN402, which contain the *Par*I-F2 fragment required for thermolabile runaway replication as well as a selectable marker gene pKN402. Capital letters indicate the *Par*I fragments of pKN402. Plasmid pCP39 lacks a 1700 bp *Par*I fragment present in pCP3, which represents part of the pKN402 *Par*I-C fragment. The ampicillin resistance gene (*bla*) and the  $\lambda$  *P<sub>L</sub>* promoter (black bars with arrow) are derived from the pPLA series of plasmids described in Fig. 7-27 (Renaut *et al.*, 1983).

quite frequently, the DNA fragment must be protected by treatment with *Hpa*II methylase prior to *Hpa*II digestion. In principle, this strategy can be used for any other gene. A disadvantage is that it does not directly allow selection of or quick screening for maximally expressing clones. In our example (Fig. 7-69) screening of a large number of transformants yielded the expression vector pCA246, which produces up to 10% of the newly synthesised protein as HBcAg after induction with 3- $\beta$ -indolylacrylic acid.

#### 7.4.2.3 The $\lambda$ *P<sub>L</sub>* System

The strong  $\lambda$  *P<sub>L</sub>* promoter has been particularly useful for the high-level expression of proteins *E. coli* cells transformed with both the runaway replication *P<sub>L</sub>* vector and the pCB57 vector. In an elegant and most efficient application it is employed in a two plasmid system (Renaut *et al.*, 1983) which also exploits the temperature-sensitive runaway repli-

cation phenomenon alluded to earlier (Section 4.1.1). One plasmid component is derived from plasmid pKN402, a 7.8 kb mini-derivative runaway replication mutant of plasmid R1c (Fig. 7-70). This plasmid contains both the temperature-sensitive replicon and the  $\lambda$  *P<sub>L</sub>* promoter; the latter lies upstream from a polylinker, which the desired gene can be inserted. Expression of the *P<sub>L</sub>* promoter from such a construct is regulated by the *cl* gene product encoded by a single chromosomal gene copy or, even better, a *cl* gene on a compatible multicopy plasmid. Such a plasmid, pCB57, is described in Section 4.1.5. It confers kanamycin resistance, and is controlled by the  $\lambda$  repressor, and is compatible with the replicon of pKN402 and its derivatives. *E. coli* cells transformed with both the runaway replication *P<sub>L</sub>* vector and the pCB57 vector contain approximately 30-50 copies of each at 28°C. At this temperature the active *c* repressor acts *in trans* to prevent any transcrip-

from the  $P_L$  promoter on the other plasmid. A shift to higher temperature (42°C) leads to two events, a ten- to twentyfold amplification of the runaway replication vector copy number, and a simultaneous derepression of the  $P_L$  promoter due to inactivation of the cI857 repressor at 42°C. This two-plasmid expression system was tested with the T4-derived DNA ligase gene, the expression of which could be induced to levels up to 25% of the total cellular protein. It is effective in many *E. coli* strains and has also proved successful for the expression of the human IGF-I protein (Buell *et al.*, 1985).

#### 7.4.2.4 Synthetic Ribosome Binding Sites

The hybrid ribosome binding sites discussed in Sections 7.4.2.1 and 7.4.2.2 are not necessarily optimal for ribosome binding, and hence for efficient translation (*cf.* also Section 7.2). These binding sites contain naturally occurring S/D sequences which frequently show a relatively low degree of homology with the sequence of the 3' end of 16S ribosomal RNA. In the *lac* system it is only four and in the *trp* leader peptide S/D sequence only three bases which show this homology at all. It was postulated (Jay *et al.*, 1981) that ribosome binding, and hence initiation of protein biosynthesis, would be much more efficient if these regions of homology could be extended. A DNA oligomer containing an S/D sequence of nine base pairs and an additional sequence,

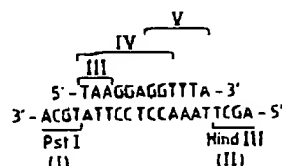


Fig. 7-71. Structure of a synthetic linker with *Pst* I (I) and *Hind* III (II) ends, coding for a stop codon (III), an S/D sequence (IV), and the GGTTTA sequence. (Jay *et al.*, 1981).

5'-GGTTTAA-3', which is important for binding ribosomal proteins (Fig. 7-71; *cf.* also Fig. 7-37; also Jay *et al.*, 1982) therefore was synthesised chemically. The entire synthetic ribosome binding site consists of two oligonucleotides of twelve and twenty bases, respectively. The left-hand 3' protruding end contains a sequence which allows ligation with a *Pst* I site (I), the right-hand 5' protruding end a *Hind* III site (II). A TAA stop codon (III) within this linker molecule is in phase with  $\beta$ -lactamase (*sec* below); in the inner part of this linker lie the S/D sequence of nine bases (IV) and the sequence GGTTTAA (V). Since the linker is asymmetrical it is more universally applicable than conventional symmetrical linkers.

As shown in the example in Fig. 7-72, this linker is positioned at a correct distance in front of the start codon of a gene to be expressed, and inserted together with this gene X into the *Pst* I site within the  $\beta$ -lactamase gene of pBR322 (*cf.* also Fig. 4.1-11). In a bacterial cell, transcription initiates at the promoter of the  $\beta$ -lactamase gene to yield a hybrid mRNA containing the  $\beta$ -lactamase component and sequences of gene X

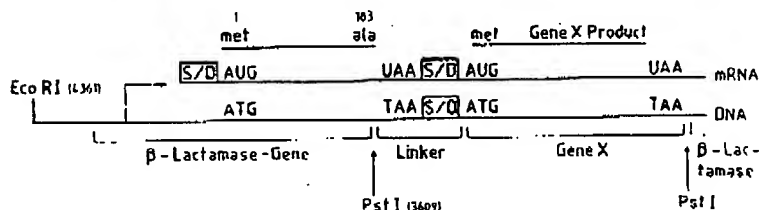


Fig. 7-72. Use of synthetic ribosomal binding sites for the construction of expression vectors. The linker carries a stop codon and a consensus S/D sequence. Although only one hybrid mRNA is transcribed, two proteins are synthesised, one of which is a fragment of  $\beta$ -lactamase with amino acids 1-183; the other is the gene X product with an N-terminal methionine residue. Numbers in parenthesis are pBR322 co-ordinates.